

Better Process Supervision with Bi-directional Rewarding Signals

Anonymous ACL submission

Abstract

Process supervision, i.e., evaluating each step, is critical for complex large language model (LLM) reasoning and test-time searching with increased inference compute. Existing approaches, represented by process reward models (PRMs), primarily focus on rewarding signals up to the current step, exhibiting a one-directional nature and lacking a mechanism to model the distance to the final target. To address this problem, we draw inspiration from the A* algorithm, which states that an effective supervisory signal should simultaneously consider the incurred cost and the estimated cost for reaching the target. Building on this key insight, we introduce **BiRM**, a novel process supervision model that not only evaluates the correctness of previous steps but also models the probability of future success. We conduct extensive experiments on mathematical reasoning tasks and demonstrate that BiRM provides more precise evaluations of LLM reasoning steps, achieving an improvement of 3.1% on Gaokao2023 over PRM under the Best-of-N sampling method. Besides, in search-based strategies, BiRM provides more comprehensive guidance and outperforms ORM by 5.0% and PRM by 3.8% respectively on MATH-500.

1 Introduction

With the rapid development of LLMs, how to supervise them has become a key research challenge, especially for complex tasks like long-term reasoning (Zelikman et al., 2022; OpenAI, 2024b; Wan et al., 2024). Previous work has explored training process supervision models to provide dense supervision on each step (Uesato et al., 2022; Lightman et al., 2024; Wang et al., 2024b), which is intuitively and practically better than outcome supervision models (Cobbe et al., 2021) that only provide sparse signals on the final answer. During test-time, process supervision models can further guide the search of LLMs or perform solution re-ranking by allocating

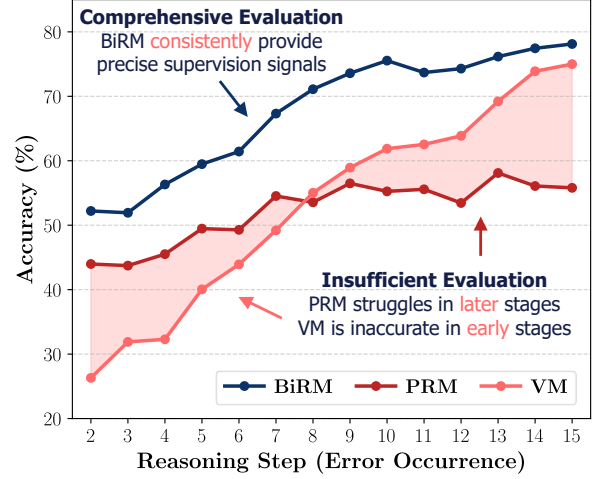


Figure 1: Error-detection accuracy across different steps, where step 1 and steps beyond 15 are truncated for better visualization. We evaluate the process reward model (PRM), value model (VM), and **BiRM** on PRMBench.

more inference compute (Snell et al., 2024; Brown et al., 2024; Wu et al., 2024).

However, existing approaches, represented by process reward models (PRMs) from OpenAI (Lightman et al., 2024), typically focus on providing one-directional reward signals on the reasoning steps that have already been generated, without consciously considering the probability of future success (Yu et al., 2024a; Zhang et al., 2025). Specifically, while they can accurately distinguish between correct and incorrect steps at the current state (i.e., backward supervision), their ability to identify which partial solution is most likely to reach the correct final answer (i.e., forward supervision) is not guaranteed, leading to sub-optimal performance in guiding effective next-step reasoning (Stroebl et al., 2024; Wang et al., 2025).

To address this challenge, we draw inspiration from the classic A* algorithm, and introduce **BiRM**, a novel process supervision model that provides bidirectional rewarding signals. Classically, the A* algorithm (Hart et al., 1968) states that an

appropriate supervisory signal should take two aspects into account: the cumulative cost up to the current step, and the estimated probability of reaching the target (Zhuang et al., 2024; Wang et al., 2024a). Motivated by this key insight, we redesign the process supervision signals, which should not only assess the correctness of steps taken so far, but also evaluate the future success probability of the partial solution. Specifically, BiRM introduces a value model (VM) head to help model the forward supervision signal (Yu et al., 2024a; Ankner et al., 2024), so that it can estimate both the correctness and success probability of a reasoning prefix/partial solution (Section 4).

To validate our motivation, we conduct a preliminary analysis on PRMBENCH (Song et al., 2025), a benchmark designed to evaluate the capability of process supervision models. We include PRM and VM as baselines, where the former estimates the correctness of partial solutions, and the latter estimates the future success probability. As shown in Figure 1, PRM performs better at detecting error steps in the early stages of reasoning, while VM performs better in the later stages. This indicates that each baseline has limitations, which aligns with the intuition we derive from the A* algorithm. In contrast, BiRM outperforms both of them in all stages, demonstrating the comprehensiveness and effectiveness of our approach.

We then perform extensive experiments on three mathematical reasoning tasks: GSM8K, MATH-500 and Gaokao2023 (Liao et al., 2024) to demonstrate the effectiveness of BiRM across different model series and search strategies. For example, BiRM trained on Qwen2.5-7B-Base achieves a 3.1% improvement on Gaokao2023 over PRM using Best-of-N sampling. Additionally, in beam search with a total sampling size of 100, BiRM further surpasses PRM by 3.8% and ORM by 5.0%.

In summary, our contributions are as follows:

- We draw inspiration from A* algorithm and propose BiRM, a novel process supervision model that provides bidirectional rewarding signals.
- We conduct extensive experiments on math reasoning tasks to demonstrate its effectiveness in solution re-ranking and trajectory searching.
- We present an in-depth analysis and demonstrate that BiRM is orthogonal to existing open-source supervision models, highlighting its robustness and generalization capabilities.

2 Related Work

2.1 Enhancing Mathematical Reasoning Capabilities of LLMs

Mathematical reasoning tasks remain a significant challenge for LLMs (OpenAI, 2024a; Snell et al., 2024). Researchers have conducted extensive studies on both train-time and test-time improvements. At train-time, supervised fine-tuning is a well-established approach. Its core idea is to construct large-scale, high-quality datasets to enhance performance (Liao et al., 2024; Yu et al., 2024b; Tong et al., 2024). On the other hand, experimental results from Openai-o1 (OpenAI, 2024b) and DeepSeek-R1 (DeepSeek-AI et al., 2025) highlight the promising potential of test-time scaling laws. Vanilla sampling methods like Best-of-N sampling (Liu et al., 2025) and search-based strategies such as beam search, A*, and MCTS (Zhuang et al., 2024; Wan et al., 2024; Zhang et al., 2024a) have all achieved remarkable performance by allocating more computational resources at test-time. In this work, we focus on improving LLM’s performance during the test-time phase.

2.2 Process Supervision Models in LLM Reasoning

LLMs can leverage an additional supervision model to achieve accurate test-time reasoning. Mainstream approaches can be divided into outcome reward models (ORMs) and process reward models (PRMs). ORMs are trained with rule-based labeled data and assign one score to the entire solution path (Cobbe et al., 2021; Yu et al., 2024a). This method achieves striking results in reasoning models like Deepseek-R1 but struggles with other tasks where the answers are highly open-ended. On the other hand, PRMs evaluate each intermediate steps in the trajectory, providing more granular reward signals (Lightman et al., 2024; Uesato et al., 2022; Zhang et al., 2025). Depending on the practical implementation, there are several variants of PRMs: (1) *Value Models* (VMs, Wang et al., 2024b; Luo et al., 2024) use Monte Carlo estimation to label steps, reducing the burden of manual annotation. The resulting labels represent the probability of future success, essentially making PRMs a type of value model. (2) *Generative Reward Models* (Zhang et al., 2024c) leverage the text generation capabilities of LLMs, providing natural language feedback, rather than traditional numerical scores.

3 Motivation

3.1 Task Formulation

Given a mathematical question q , a large language model π generates a sequence of reasoning steps to solve the problem. The complete reasoning trajectory, i.e., chain-of-thought (Wei et al., 2022), can be denoted as $\tau = \{s_1, s_2, \dots, s_m\}$, where s_i represents the i -th step and m is the number of total reasoning steps.

3.2 The Limitations of PRMs

PRMs are typically trained to assign a numerical score to each intermediate reasoning step, evaluating their correctness. For a partial trajectory $\tau^{[1:t]} = \{s_1, s_2, \dots, s_t\}$, PRM can provide an reward score for step s_i :

$$r(s_i, q) = p(s_i \text{ is correct} \mid q), \quad (1)$$

where $r(\cdot)$ represents the process-based reward function provided by PRM. Further, the correctness of the partial trajectory $\tau^{[1:t]}$ can be expressed as the accumulative correctness reward of all intermediate steps, following Lightman et al. (2024):

$$\begin{aligned} \mathcal{R}(\tau^{[1:t]}, q) &= p([s_1, s_2, \dots, s_t] \text{ is correct} \mid q) \\ &= \prod_{i=1}^t p(s_i \text{ is correct} \mid q) = \prod_{i=1}^t r(s_i, q). \end{aligned}$$

This equation highlights the one-directional scoring nature of PRMs, which evaluate whether the sampled trajectory $\{s_1, s_2, \dots, s_t\}$ is correct given the problem q . Instead, for the potential future paths $\{s_{t+1}, s_{t+2}, \dots, s_m\}$ starting from the current state s_t , PRMs lack the capability to provide effective guidance, as Figure 2 illustrates.

3.3 Inspiration from the A* Search Algorithm

To address this limitation, we draw inspiration from the A* algorithm. Originally, A* is a heuristic graph search algorithm designed to find the optimal path (Hart et al., 1968). The key insight from A* is that a good supervision signal should simultaneously consider two aspects: the accumulative cost $g(n)$ up to the current step and the future cost $h(n)$ to the target. The final value of a step is given by $f(n) = g(n) + h(n)$.

In the context of LLM mathematical reasoning, we argue that a good supervision signal should not only consider the correctness of previous steps (i.e., backward supervision) but also model the probability of future success (i.e., forward supervision). On

the one hand, PRM can naturally function as $g(\cdot)$. In other words, PRM can use its one-directional scoring ability to provide rewards for the partial solution up to the current step s_t :

$$g(s_t) = \text{Agg}(r(s_1), r(s_2), \dots, r(s_t)) = \mathcal{R}(\tau^{[1:t]}),$$

where $\text{Agg} \in \{\prod, \min, \max, \text{avg}\}$ stands for an aggregation function to summarize the accumulative rewards of all steps from s_1 to s_t .

On the other hand, to heuristically model the probability of reaching the correct final answer, we seek to utilize a value model (VM) to play the role of $h(\cdot)$. For the partial solution $\tau^{[1:t]}$, a forward-looking VM can provide a reliable probability estimation:

$$\begin{aligned} h(s_t) &= \mathcal{V}(\tau^{[1:t]}, q) \\ &= \mathbb{E}_{\hat{a} \sim \pi(\cdot \mid \tau^{[1:t]}, q)} [p(\hat{a} \text{ is correct} \mid q)]. \end{aligned} \quad (2)$$

Here, \hat{a} represents the final answer predicted by the LLM π , and $\mathcal{V}(\cdot)$ denotes the estimation of VM for whether the partial trajectory can reach the correct answer. In practical implementations, the VM and PRM share the same model architecture, but differ in the meaning of training labels, which fundamentally trains the VM as a reliable predictive estimator. We will discuss more details in Section 4.2. Finally, the complete value function can be expressed as:

$$f(s_t) = g(s_t) + \beta \cdot h(s_t), \quad (3)$$

where the coefficient β balances the importance of the $g(s_t)$ and $h(s_t)$ terms. When a step s_i has a higher $f(s_i)$ value, it indicates that this step is more promising among multiple candidates, thus contributing to more effective next-step reasoning.

4 BiRM, a Bidirectional Process Supervision Model

4.1 Training Methodology

For a query q from the training question set \mathcal{Q} , we first sample N solutions from the generator π . Then, we annotate each intermediate step of these solutions, i.e., annotating step-level labels. The resulting dataset \mathcal{D} for query q can be formalized as $\mathcal{D}_q = \{\tau_i, \{y_i^1, y_i^2, \dots, y_i^j, \dots\}\}_{i=1}^N$, where τ_i denotes the i -th sampled trajectory, and y_i^j represents the step label for the j -th step in the i -th solution. We will introduce more annotation strategies in Section 4.2.

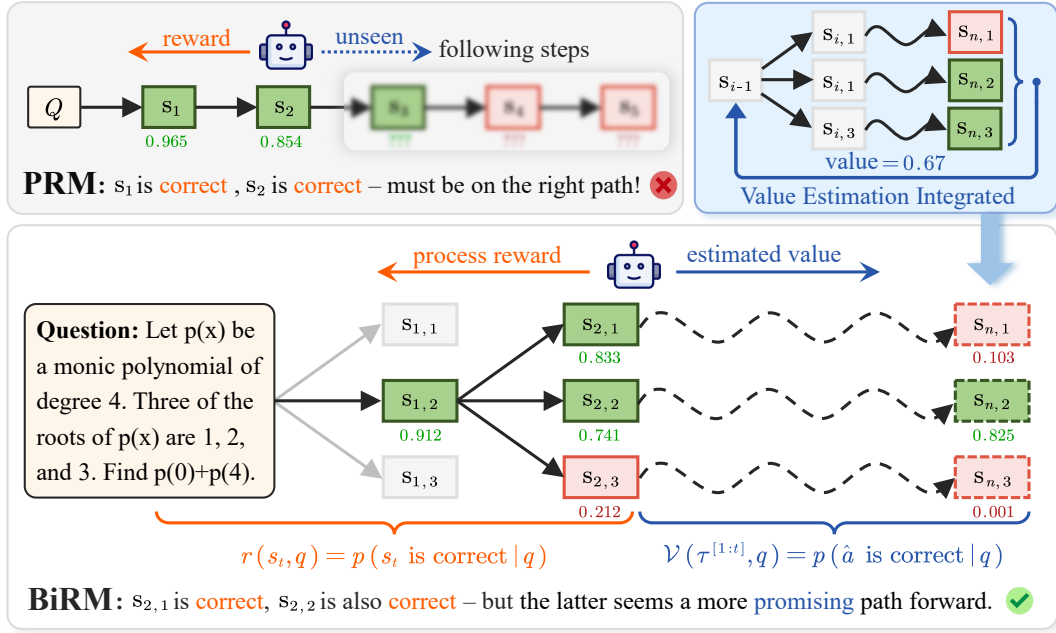


Figure 2: An example of our proposed **BiRM** compared with traditional Process Reward Models (PRMs). Given a question q , PRMs only consider the accumulated rewards up to the current step. In contrast, BiRM takes into account two aspects: the correctness rewards received so far and the probability of reaching correct final answers.

Following Yu et al. (2024a), we implement the vanilla PRM by adding a linear layer for reward prediction after the last hidden layer of the LLM. We also retain the original language modeling head. Formally, a vanilla PRM $\mathcal{R}(\theta, \phi_R)$ is parameterized by base model parameters θ and reward head parameters ϕ_R . The training objective of PRM is to minimize the mean squared error (MSE) loss between the predicted reward scores and the binary step-level reward labels. Thus, we have:

$$\mathcal{L}_{\text{PRM}}(\theta, \phi_R) = \frac{1}{|\mathcal{Q}|} \sum_{q \in \mathcal{Q}} \left[\mathbb{E}_{\tau \sim \pi(\cdot|q)} \sum_{t=1}^m (\hat{r}_{\theta, \phi_R}(s_t, q) - r^t)^2 \right],$$

where $\hat{r}(s_t, q)$ represents the predicted reward score for the t -th step (Equation 1), and r^t denotes the ground truth step label. m represents the total number of steps in solution τ .

Furthermore, to alleviate the one-directional limitation of PRM, we introduce an additional value head to guide process supervision. Specifically, BiRM $\mathcal{M}(\theta, \phi_R, \phi_V)$ is parameterized by three components: θ represents the base model parameters, ϕ_R represents the reward head, and ϕ_V corresponds to the value head. The overall training objective of BiRM is to jointly minimize the discrepancy between the predicted reward score and the reward label, as well as between the value score

and the value label. Similar to the vanilla PRM, we employ MSE loss for the value head:

$$\mathcal{L}_{\text{VM}}(\theta, \phi_V) = \frac{1}{|\mathcal{Q}|} \sum_{q \in \mathcal{Q}} \left[\mathbb{E}_{\tau \sim \pi(\cdot|q)} \sum_{t=1}^m (\hat{\mathcal{M}}_{\theta, \phi_V}(\tau^{[1:t]}, q) - v^t)^2 \right],$$

where $\hat{\mathcal{M}}_{\theta, \phi_V}$ represents the estimated success probability for the partial solution $\tau^{[1:t]}$ (Equation 2), and v^t denotes the value label for s_t .

In this way, the optimized BiRM considers not only the actual accumulative rewards obtained so far, but also the potential of reaching correct final answers (Figure 2). The complete loss function for BiRM can be defined as:

$$\mathcal{L}_{\text{BiRM}}(\theta, \phi_R, \phi_V) = \mathcal{L}_{\text{PRM}}(\theta, \phi_R) + c \cdot \mathcal{L}_{\text{VM}}(\theta, \phi_V). \quad (4)$$

We use a coefficient c to balance the importance of the reward term \mathcal{L}_{PRM} and the value term \mathcal{L}_{VM} .

4.2 Step Label Annotation Strategies

In this section, we discuss our annotation strategies for two kinds of BiRM training labels.

Reward Labels. Reward labels are defined as the correctness of each current step, represented as a binary label. We use the MetaMath dataset (Yu et al., 2024b) as our training data. We first perform

supervised fine-tuning on the base model to obtain the generators. Then, we sample 15 rollouts for each query and use Deepseek-V3 (DeepSeek-AI et al., 2024) to annotate the correctness of each step. Detailed annotation procedures and prompts are provided in Appendix B.2.

Value Labels. A key challenge in implementing the value head is to accurately estimate value labels for the partial solution $\tau^{[1:t]}$. We employ multiple strategies to address this problem.

MC-based estimation (Wang et al., 2024b) is a widely used method for automated labeling, which can be categorized into soft-label and hard-label annotations. Specifically, we sample N rollouts from an intermediate step in the trajectory. If M of them are correct, the soft-label for the current step can be defined as: $\text{label}(s_t) = \frac{M}{N}$. In contrast, the hard-label method suggests that if any of the rollouts reaches the target, then $\text{label}(s_t) = 1$.

The essence of Monte Carlo estimation is to assess the potential of reaching correct final answer from the current step and assign this probability to the step label. Thus, for estimating a partial trajectory, we can formally express it as:

$$\mathcal{V}(\tau^{[1:t]}, q) \approx \frac{1}{N} \sum_{i=1}^N \mathbb{I}(\hat{a}_i \text{ is correct} \mid \tau^{[1:t]}, q).$$

As the number of rollouts N increases, the estimated value label becomes more accurate. Following Wang et al. (2024b), we sample 8 solutions for each intermediate step and analyze the effectiveness of both soft-label and hard-label approaches.

Outcome-supervised estimation (Yu et al., 2024a) states that using the outcome label alone is sufficient to provide probability estimation for each reasoning steps. The underlying idea is that during the training phase, we can replicate the final answer’s correctness label across all intermediate steps. The resulting value model implicitly learns to foresee the future, predicting potential final outcome (i.e. value) for partial solutions. Compared to MC estimation, outcome-supervised estimation has higher data efficiency, but the shortcoming is that the automatically learned estimation in this way is less accurate.

5 Experiments

5.1 Experimental Setup

Tasks. We conduct experiments using three widely used math reasoning datasets: GSM8K

(Cobbe et al., 2021), MATH-500 (Lightman et al., 2024), and an out-of-domain (OOD) dataset Gaokao2023 (Liao et al., 2024) to evaluate the generalization ability of BiRM. Besides, we test our method on three base models across different model sizes and families: Qwen2.5-3B-Base (Yang et al., 2024), Qwen2.5-7B-Base (Yang et al., 2024), and Llama3.1-8B-Base (Dubey et al., 2024).

Baselines. To verify the effectiveness of BiRM, we consider a wide range of baselines, including the outcome reward model (ORM, Cobbe et al., 2021), process reward model (PRM, Lightman et al., 2024) and two variants of PRM: Math-Shepherd (Wang et al., 2024b) and ER-PRM (Zhang et al., 2024b). Additionally, we include the results of greedy decoding and rule-based approaches, i.e. Majority Voting. We present more details in Appendix A.1.

Implementation Details. In the SFT phase, we train our generators on the MATH subset of the MetaMath dataset (Yu et al., 2024b) for two epochs, with a learning rate set to 1×10^{-5} . The global batch size is set to 256. In the training phase, we use 225,000 sampled solutions to train BiRM for one epoch based on the generator checkpoint with a learning rate of 5×10^{-6} . More details are provided in Appendix A.2.

Evaluation Metrics. We conduct a comprehensive evaluation of BiRM, considering both vanilla sampling and search strategies. Best-of-N (BoN) sampling is a commonly used evaluation metric for PRMs. It requires the model to score N candidate solutions, with the highest-scoring solution selected as the final outcome. We also conduct beam search experiments to verify that BiRM can provide more promising guidance for LLM reasoning. In practice, BiRM follows Equation 3, estimating both rewards and values to calculate final scores. A detailed description is provided in Appendix A.3.

5.2 Main Results

BiRM exhibits more comprehensive and superior evaluations in BoN sampling. Table 1 presents a comparison of BoN accuracy across different supervision models on GSM8K, MATH-500, and the out-of-domain Gaokao2023 dataset. Our observations are as follows: (1) BiRM consistently outperforms vanilla ORM, PRM, and their variants on both GSM8K and MATH-500. For instance, BiRM trained on Llama3.1-8B outperforms

Models	Methods	Avg.	GSM8K			MATH-500			Gaokao2023		
			@128	@256	@512	@128	@256	@512	@128	@256	@512
Qwen2.5-3B	Greedy	46.8	—	73.1	—	—	40.2	—	—	27.0	—
	Majority Vote	58.1	85.1	85.0	85.3	52.5	53.0	53.8	35.8	36.3	36.1
	ORM	58.9	88.1	88.1	88.1	52.1	51.8	52.2	37.2	37.0	35.8
	PRM	59.9	88.5	88.3	88.0	54.6	54.1	54.2	37.3	37.1	37.2
	ER-PRM	58.8	88.0	88.0	87.7	52.6	52.3	52.0	36.2	36.3	35.8
	Math-Shepherd	59.0	87.3	87.2	87.0	53.2	53.4	53.8	36.6	36.4	36.1
	BiRM	61.0	88.4	88.6	88.9	55.9	56.1	57.4	36.9	37.8	38.7
Qwen2.5-7B	Greedy	52.3	—	78.5	—	—	45.0	—	—	33.5	—
	Majority Vote	63.6	88.1	88.0	87.8	57.3	57.5	57.6	45.5	45.4	45.2
	ORM	64.7	92.0	91.6	91.3	59.6	59.9	59.4	43.6	43.5	41.3
	PRM	66.3	92.7	92.8	92.9	60.3	60.1	58.4	45.8	46.2	47.3
	ER-PRM	66.2	92.2	92.1	92.2	59.7	59.2	59.0	47.0	47.2	47.3
	Math-Shepherd	66.3	92.1	92.2	91.7	60.3	60.2	60.4	46.4	47.0	46.5
	BiRM	68.3	93.1	93.3	93.2	62.4	62.3	63.4	47.7	49.1	50.4
Llama3.1-8B	Greedy	34.7	—	55.7	—	—	31.2	—	—	17.1	—
	Majority Vote	46.4	72.1	72.0	72.3	39.2	40.2	41.1	26.5	27.2	27.1
	ORM	50.3	84.1	84.5	85.0	41.5	40.9	40.8	25.4	25.2	24.9
	PRM	51.5	84.1	84.8	85.2	42.5	42.2	41.8	28.2	27.7	27.3
	ER-PRM	50.6	84.8	85.3	85.8	41.3	41.0	40.2	25.7	26.1	24.9
	Math-Shepherd	51.3	84.4	84.9	85.3	42.7	42.9	43.6	25.8	25.8	26.2
	BiRM	54.1	86.1	87.2	87.8	45.4	45.4	45.6	29.4	30.0	29.6

Table 1: Performance of Best-of-N sampling on GSM8K, MATH-500 and Gaokao2023 with three base models. The accuracy of the BoN solution is utilized as the evaluation metric. The results are reported as the average accuracy across five random seeds. @128, @256, and @512 denote the accuracy with Best-of-128, Best-of-256, and Best-of-512 sampling, respectively. The results of greedy decoding are independent of N and are listed for comparison purposes. The best results are marked in **bold**.

PRM on GSM8K by 2.6%, while BiRM based on Qwen2.5-7B achieves an additional 5.0% improvement on MATH-500. (2) BiRM exhibits better generalization ability. Since supervision models are trained solely on the query sets from GSM8K and MATH, Gaokao2023 serves as an out-of-domain (OOD) test set. BiRM-Qwen2.5-7B surpasses the finely labeled Math-Shepherd by 3.9%. In contrast, other supervision methods show fluctuating performance across different base models. (3) As N increases, some supervision methods fail to provide consistent supervision. For example, ORM trained on Qwen2.5-3B shows a decrease on Gaokao2023 from 37.2% to 35.8%. In contrast, BiRM maintains a continuous increase in accuracy. We provide more detailed discussions in Section 6.1.

BiRM demonstrates more meaningful and promising guidance in search-based strategies. To fully demonstrate the superiority of BiRM’s bidirectional supervision capability, we conduct further experiments under search-based strategies. We run step-level beam search and choose vanilla ORM and PRM as baselines. The detailed algo-

rithm is provided in Appendix A.3. From Table 2, we can conclude that: (1) BiRM achieves the highest accuracy in most cases. For example, on GSM8K, Qwen2.5-7B-BiRM achieves an accuracy of 89.4 at $K = 8$, which is a notable improvement over PRM’s 88.1%. (2) As beam size increases, BiRM’s performance continues to improve. On the Llama3.1-8B base model, BiRM outperforms ORM by 2.8% at $K = 20$ and achieves a notable 5.0% improvement at $K = 100$ in MATH-500 dataset. These results emphasize the valuable bidirectional supervision signals provided by BiRM, which significantly contributes to guiding the LLM toward more successful and promising final answers in solution searching.

6 Analysis and Discussions

6.1 Scaling Decline in BoN sampling

We conduct a further analysis of the scaling decline phenomenon in our main results. The complete BoN accuracy curve, shown in Figure 3, is plotted for N ranging from 1 to 512. As N increases, we observe that BiRM shows a consistent improve-

Models	# Total Size	GSM8K			MATH-500			Gaokao2023		
		ORM	PRM	BiRM	ORM	PRM	BiRM	ORM	PRM	BiRM
Qwen2.5-3B	$K = 4$	83.0	82.1	82.8	48.6	49.3	50.1	35.6	34.9	36.1
	$K = 8$	84.6	83.9	85.1	50.1	50.9	52.5	36.1	37.9	37.9
	$K = 20$	86.7	85.7	86.9	53.0	54.3	55.0	37.7	38.4	39.1
	$K = 100$	87.5	85.9	87.6	53.0	53.9	55.1	38.1	37.9	39.0
Qwen2.5-7B	$K = 4$	86.2	86.5	87.0	55.7	55.8	57.1	42.8	44.0	44.5
	$K = 8$	88.6	88.1	89.4	58.3	59.1	60.1	44.2	45.6	46.8
	$K = 20$	90.4	89.2	90.6	59.1	61.5	62.3	45.5	48.1	48.4
	$K = 100$	91.2	88.4	91.7	60.1	60.7	62.5	46.8	48.3	50.0
Llama3.1-8B	$K = 4$	72.8	71.7	72.9	38.5	39.9	40.7	23.9	25.1	25.4
	$K = 8$	77.4	75.9	78.3	40.2	40.1	43.3	25.6	26.6	27.5
	$K = 20$	81.4	79.2	81.7	41.5	42.1	44.3	27.0	28.6	29.2
	$K = 100$	82.7	80.3	85.4	41.1	42.3	46.1	26.2	29.6	30.7

Table 2: Performance of beam search on GSM8K, MATH-500 and Gaokao2023 with three base models. “# Total Size” stands for total sampling size K in beam search and we report the best performance among all beam sizes. The results are reported as the average accuracy across three random seeds. The best results are marked in **bold**.

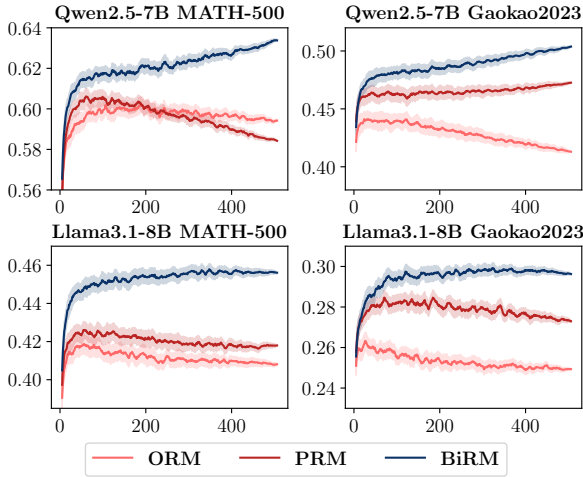


Figure 3: Scaling decline phenomenon in Best-of-N sampling. We present the BoN accuracy results across five random seeds. For better visualization, we apply the moving average with a window size of 10.

ment. In contrast, the post-verification accuracy of vanilla ORM and PRM plateaus and even declines, which contradicts our intuition learned from the test-time scaling laws (Snell et al., 2024).

We attribute this decline to verifier failures. Imperfect verifiers misrank candidates, erroneously classifying positive samples as negative. As the sample size increases, this misjudgment becomes more pronounced. Traditional PRMs exhibit a one-directional scoring nature, limiting their ability to evaluate candidates from a comprehensive perspective. In contrast, BiRM estimates both rewards and values, providing more reliable supervision signals.

Models	Methods	MATH-500		Gaokao2023	
		@128	@512	@128	@512
Qwen2.5-7B	+ Outcome	61.8	61.1	46.8	49.4
	+ MS. (Hard)	62.1	62.8	47.3	49.7
	+ MS. (Soft)	62.4	63.4	47.7	50.4
Llama3.1-8B	+ Outcome	44.9	44.2	29.0	29.6
	+ MS. (Hard)	45.1	45.4	29.2	29.4
	+ MS. (Soft)	45.4	45.6	29.4	29.6

Table 3: Different value label annotation strategies for BiRM. “Outcome” stands for Outcome-supervised estimation. “MS. (Hard)” and “MS. (Soft)” represents Math-Shepherd hard and soft estimation respectively.

6.2 Annotation Strategies for Value Labels

As discussed in Section 4.2, we explore various strategies for annotating precise value labels. We aim to demonstrate that our method has good orthogonality with existing annotation strategies.

Table 3 presents the accuracy of BiRM in BoN sampling under different strategies. We can conclude that: (1) More accurate annotations lead to greater improvements. The mash-shepherd soft estimation, which uses the potential success probability of intermediate steps as explicit labels, offers the finest granularity and achieves the best performance. In contrast, outcome-supervised estimation, which relies on outcome labels for implicit learning, exhibits greater variability. (2) Even the weakest method, outcome-supervised estimation, shows a notable improvement over PRM. This highlights the flexibility and applicability of BiRM.

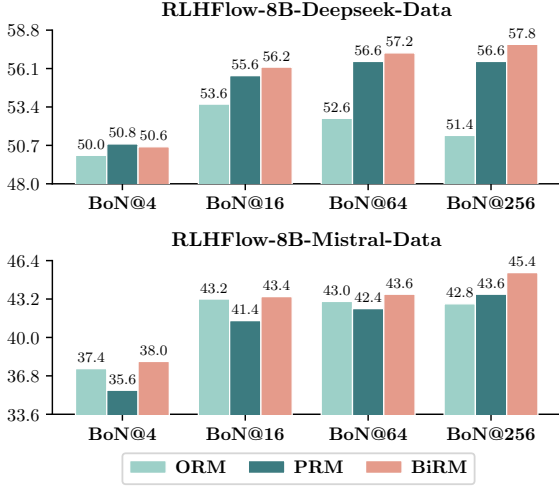


Figure 4: Performance comparison of ORM, PRM and BiRM under BoN sampling. The base models are open-source RLHFlow-8B-Deepseek-Data and RLHFlow-8B-Mistral-Data (Xiong et al., 2024). We follow Equation 3 to calculate the BiRM score at test-time.

6.3 Orthogonality to Existing PRMs

To further demonstrate the generalization ability of our method, we conduct experiments using several existing open-source reward models. We select ORMs and PRMs trained by RLHFlow (Xiong et al., 2024; Dong et al., 2024) as baselines and reuse the N sampled solutions they provided for testing. Then we follow Equation 3 to calculate the BiRM scores for BoN sampling.

Experiment results in Figure 4 clearly reveal that BiRM consistently outperforms both ORM and PRM across different values of N , maintaining a consistent upward trend. Furthermore, this trend expands at larger sampling sizes, where BiRM maintains its lead, reaching an accuracy of 57.8% at BoN@256, compared to PRM’s 56.6% and ORM’s 51.4%, respectively. These findings indicate the reliability and generalization of BiRM when using existing open-source reward models.

6.4 Query Scaling or Response Scaling

We also explore a key issue in training supervision models: which matters more, query scaling or response scaling?

We first fix the number of queries and use the original GSM8K and MATH datasets, which contain approximately 15,000 queries. We then test BiRM’s performance with response sizes of 8, 15, and 30. The results in Table 4 reveal that BiRM performs best when the response size is 15 on both datasets. The possible reason is that when the

# Query	# Resp.	MATH-500		Gaokao2023	
		@128	@512	@128	@512
15,000	$\times 30$	61.3	61.6	47.3	48.3
	$\times 15$	62.0	63.0	46.8	49.4
	$\times 8$	61.3	61.2	46.4	46.8
7,500	$\times 15$	59.0	58.8	45.4	44.7
3,750	$\times 30$	57.9	58.2	43.4	42.8

Table 4: Training data scaling for queries and responses. The base model is Qwen2.5-7B and we use outcome-supervised estimation for simplicity.

number of responses is too low, BiRM cannot learn sufficient and diverse supervision signals. On the other hand, the model struggles with overly similar data patterns per query when # Resp. = 30, leading to overfitting.

Furthermore, we control the total size of the training dataset. Specifically, we conduct experiments with three following settings: $15,000 \times 8$, $7,500 \times 15$, and $3,750 \times 30$. The results demonstrate that BiRM performs best with the $15,000 \times 8$ configuration. Additionally, we observe that models with fewer queries go through more severe degradation when facing OOD test sets. In the MATH-500 experiments, the gap between the $7,500 \times 15$ and $3,750 \times 30$ settings ranges from 0.6% to 1.1%, but this gap significantly widens to 2.0% on the Gaokao2023 benchmark. To sum up, we believe that maintaining an appropriate response size while scaling the number of queries is critical to training process supervision models. We hope this provides valuable insights to the community.

7 Conclusion

In this work, we introduce BiRM, a novel process supervision model for large language models (LLMs), inspired by the A* algorithm. BiRM provides bidirectional supervision signals, evaluating both the correctness of reasoning steps taken so far and the probability of reaching correct answers in the future. Our extensive experiments demonstrate the effectiveness of BiRM across various mathematical reasoning tasks, outperforming existing supervision models like ORM and PRM. Through detailed analysis, we highlight the strengths of BiRM in guiding the search process and improving solution re-ranking. We hope that our approach contributes valuable insights to the field of process supervision and opens avenues for future research in enhancing LLM-based reasoning.

Limitations

Our work has some limitations, which we leave for future work to address: (1) High computational cost in test-time searching. In order to improve the performance of LLMs at test-time, we employ vanilla sampling and search-based strategies for solution searching. However, this process requires a significant amount of computational resources. In our work, we use vLLM (Kwon et al., 2023) to alleviate this limitation. Besides, we also observe that search-based strategies sometimes perform worse than repeated sampling due to verifier failures (Yu et al., 2025), even under the same computational budget. We will explore this problem in the future. (2) Generalization across different data patterns and base models. In our experiments, we train our generators and supervision models based on the same base models, ensuring the same data patterns. However, in practical scenarios, an optimal supervision model should be independent of the data pattern and capable of supervising different kinds of reasoning paths. We hope our work provides insights to the community and contributes to the development of more robust and generalized supervision models.

References

- Zachary Ankner, Mansheej Paul, Brandon Cui, Jonathan D. Chang, and Prithviraj Ammanabrolu. 2024. [Critique-out-loud reward models](#). *CoRR*, abs/2408.11791.
- Bradley C. A. Brown, Jordan Juravsky, Ryan Saul Ehrlich, Ronald Clark, Quoc V. Le, Christopher Ré, and Azalia Mirhoseini. 2024. [Large language mon-keys: Scaling inference compute with repeated sampling](#). *CoRR*, abs/2407.21787.
- Karl Cobbe, Vineet Kosaraju, Mohammad Bavarian, Mark Chen, Heewoo Jun, Lukasz Kaiser, Matthias Plappert, Jerry Tworek, Jacob Hilton, Reiichiro Nakano, Christopher Hesse, and John Schulman. 2021. [Training verifiers to solve math word problems](#). *CoRR*, abs/2110.14168.
- DeepSeek-AI, Daya Guo, Dejian Yang, Haowei Zhang, Junxiao Song, Ruoyu Zhang, Runxin Xu, Qihao Zhu, Shirong Ma, Peiyi Wang, Xiao Bi, Xiaokang Zhang, Xingkai Yu, Yu Wu, Z. F. Wu, Zhibin Gou, Zhihong Shao, Zhuoshu Li, Ziyi Gao, Aixin Liu, Bing Xue, Bingxuan Wang, Bochao Wu, Bei Feng, Chengda Lu, Chenggang Zhao, Chengqi Deng, Chenyu Zhang, Chong Ruan, Damai Dai, Deli Chen, Dongjie Ji, Erhang Li, Fangyun Lin, Fucong Dai, Fuli Luo, Guangbo Hao, Guanting Chen, Guowei Li, H. Zhang, Han Bao, Hanwei Xu, Haocheng Wang, Honghui Ding, Huajian Xin, Huazuo Gao, Hui Li, Hui Qu, J. L. Cai, Jian Liang, Jianzhong Guo, Jiaqi Ni, Jiashi Li, Jiawei Wang, Jin Chen, Jingchang Chen, Jingyang Yuan, Junjie Qiu, Junlong Li, J. L. Cai, Jiaqi Ni, Jian Liang, Jin Chen, Kai Dong, Kai Hu, Kaige Gao, Kang Guan, Kexin Huang, Kuai Yu, Lean Wang, Lecong Zhang, Liang Zhao, Litong Wang, Liyue Zhang, Lei Xu, Leyi Xia, Mingchuan Zhang, Minghua Zhang, Minghui Tang, Meng Li, Miaojun Wang, Mingming Li, Ning Tian, Panpan Huang, Peng Zhang, Qiancheng Wang, Qihao Zhu, Qinyu Chen, Qiushi Du, R. J. Chen, R. L. Jin, Ruiqi Ge, Ruisong Zhang, Ruizhe Pan, Runji Wang, Runxin Xu, Ruoyu Zhang, Ruyi Chen, S. S. Li, Shanghao
- Ding, Huajian Xin, Huazuo Gao, Hui Qu, Hui Li, Jianzhong Guo, Jiashi Li, Jiawei Wang, Jingchang Chen, Jingyang Yuan, Junjie Qiu, Junlong Li, J. L. Cai, Jiaqi Ni, Jian Liang, Jin Chen, Kai Dong, Kai Hu, Kaige Gao, Kang Guan, Kexin Huang, Kuai Yu, Lean Wang, Lecong Zhang, Liang Zhao, Litong Wang, Liyue Zhang, Lei Xu, Leyi Xia, Mingchuan Zhang, Minghua Zhang, Minghui Tang, Meng Li, Miaojun Wang, Mingming Li, Ning Tian, Panpan Huang, Peng Zhang, Qiancheng Wang, Qihao Zhu, Qinyu Chen, Qiushi Du, Ruiqi Ge, Ruisong Zhang, Ruizhe Pan, Runji Wang, Runxin Xu, Ruoyu Zhang, Ruyi Chen, S. S. Li, Shanghao

653	Lu, Shangyan Zhou, Shanhuang Chen, Shaoqing Wu,	<i>ACL 2024, Bangkok, Thailand and virtual meeting,</i>	713
654	Shengfeng Ye, Shengfeng Ye, Shirong Ma, Shiyu	<i>August 11-16, 2024, pages 905–924. Association for</i>	714
655	Wang, Shuang Zhou, Shuiping Yu, Shunfeng Zhou,	<i>Computational Linguistics.</i>	715
656	Shuting Pan, T. Wang, Tao Yun, Tian Pei, Tianyu Sun,		
657	W. L. Xiao, and Wangding Zeng. 2024. Deepseek-v3		
658	technical report . <i>CoRR</i> , abs/2412.19437.		
659	Hanze Dong, Wei Xiong, Bo Pang, Haoxiang Wang,	Hunter Lightman, Vineet Kosaraju, Yuri Burda, Harri-	716
660	Han Zhao, Yingbo Zhou, Nan Jiang, Doyen Sahoo,	son Edwards, Bowen Baker, Teddy Lee, Jan Leike,	717
661	Caiming Xiong, and Tong Zhang. 2024. Rlhf work-	John Schulman, Ilya Sutskever, and Karl Cobbe.	718
662	flow: From reward modeling to online rlhf . <i>Preprint</i> ,	2024. Let’s verify step by step . In <i>The Twelfth In-</i>	719
663	arXiv:2405.07863.	<i>ternational Conference on Learning Representations,</i>	720
664		<i>ICLR 2024, Vienna, Austria, May 7-11, 2024. Open-</i>	721
665	Abhimanyu Dubey, Abhinav Jauhri, Abhinav Pandey,	<i>Review.net.</i>	722
666	Abhishek Kadian, Ahmad Al-Dahle, Aiesha Letman,		
667	Akhil Mathur, Alan Schelten, Amy Yang, Angela	Yantao Liu, Zijun Yao, Rui Min, Yixin Cao, Lei Hou,	723
668	Fan, Anirudh Goyal, Anthony Hartshorn, Aobo Yang,	and Juanzi Li. 2025. Pairwise rm: Perform best-of-n	724
669	Archi Mitra, Archie Sravankumar, Artem Korenev,	sampling with knockout tournament. <i>arXiv preprint</i>	725
670	Arthur Hinsvark, Arun Rao, Aston Zhang, Aurélien	<i>arXiv:2501.13007</i> .	726
671	Rodriguez, Austen Gregerson, Ava Spataru, Bap-		
672	tiste Rozière, Bethany Biron, Binh Tang, Bobbie	Liangchen Luo, Yinxiao Liu, Rosanne Liu, Samrat	727
673	Chern, Charlotte Caucheteux, Chaya Nayak, Chloe	Phatale, Harsh Lara, Yunxuan Li, Lei Shu, Yun	728
674	Bi, Chris Marra, Chris McConnell, Christian Keller,	Zhu, Lei Meng, Jiao Sun, and Abhinav Rastogi.	729
675	Christophe Touret, Chunyang Wu, Corinne Wong,	2024. Improve mathematical reasoning in language	730
676	Cristian Canton Ferrer, Cyrus Nikolaidis, Damien Al-	models by automated process supervision . <i>CoRR</i> ,	731
677	lonsius, Daniel Song, Danielle Pintz, Danny Livshits,	abs/2406.06592.	732
678	David Esiobu, Dhruv Choudhary, Dhruv Mahajan,		
679	Diego Garcia-Olano, Diego Perino, Dieuwke Hupkes,	OpenAI. 2024a. GPT-4o .	733
680	Egor Lakomkin, Ehab AlBadawy, Elina Lobanova,		
681	Emily Dinan, Eric Michael Smith, Filip Radenovic,	OpenAI. 2024b. Introducing openai o1 .	734
682	Frank Zhang, Gabriel Synnaeve, Gabrielle Lee, Georg-		
683	ia Lewis Anderson, Graeme Nail, Grégoire Mialon,	Charlie Snell, Jaehoon Lee, Kelvin Xu, and Aviral Ku-	735
684	Guan Pang, Guillem Cucurell, Hailey Nguyen, Han-	mar. 2024. Scaling LLM test-time compute optimally	736
685	nah Korevaar, Hu Xu, Hugo Touvron, Iliyan Zarov,	can be more effective than scaling model parameters .	737
686	Imanol Arrieta Ibarra, Isabel M. Kloumann, Ishan	<i>CoRR</i> , abs/2408.03314.	738
687	Misra, Ivan Evtimov, Jade Copet, Jaewon Lee, Jan		
688	Geffert, Jana Vranes, Jason Park, Jay Mahadeokar,	Mingyang Song, Zhaochen Su, Xiaoye Qu, Jiawei Zhou,	739
689	Jeet Shah, Jelmer van der Linde, Jennifer Billock,	and Yu Cheng. 2025. Prmbench: A fine-grained	740
690	Jenny Hong, Jenya Lee, Jeremy Fu, Jianfeng Chi,	and challenging benchmark for process-level reward	741
691	Jianyu Huang, Jiawen Liu, Jie Wang, Jiecao Yu,	models. <i>arXiv preprint arXiv:2501.03124</i> .	742
692	Joanna Bitton, Joe Spisak, Jongsoo Park, Joseph		
693	Rocca, Joshua Johnstun, Joshua Saxe, Junteng Jia,	Benedikt Stroebl, Sayash Kapoor, and Arvind	743
694	Kalyan Vasuden Alwala, Kartikeya Upasani, Kate	Narayanan. 2024. Inference scaling laws: The limits	744
695	Plawiak, Ke Li, Kenneth Heafield, Kevin Stone, and	of llm resampling with imperfect verifiers. <i>arXiv</i>	745
696	et al. 2024. The llama 3 herd of models . <i>CoRR</i> ,	<i>preprint arXiv:2411.17501</i> .	746
697	abs/2407.21783.		
698		Yuxuan Tong, Xiwen Zhang, Rui Wang, Ruidong Wu,	747
699	Peter E. Hart, Nils J. Nilsson, and Bertram Raphael.	and Junxian He. 2024. Dart-math: Difficulty-aware	748
700	1968. A formal basis for the heuristic determina-	rejection tuning for mathematical problem-solving .	749
701	tion of minimum cost paths . <i>IEEE Trans. Syst. Sci.</i>	In <i>Advances in Neural Information Processing Sys-</i>	750
702	<i>Cybern.</i> , 4(2):100–107.	<i>tems 38: Annual Conference on Neural Information</i>	751
703		<i>Processing Systems 2024, NeurIPS 2024, Vancouver,</i>	752
704	Woosuk Kwon, Zhuohan Li, Siyuan Zhuang, Ying	<i>BC, Canada, December 10 - 15, 2024.</i>	753
705	Sheng, Lianmin Zheng, Cody Hao Yu, Joseph Gonza-		
706	lez, Hao Zhang, and Ion Stoica. 2023. Efficient mem-	Jonathan Uesato, Nate Kushman, Ramana Kumar,	754
707	ory management for large language model serving	H. Francis Song, Noah Y. Siegel, Lisa Wang, Anto-	755
708	with pagedattention . In <i>Proceedings of the 29th Sym-</i>	tonia Creswell, Geoffrey Irving, and Irina Higgins.	756
709	<i>posium on Operating Systems Principles, SOSP 2023,</i>	2022. Solving math word problems with process- and	757
710	<i>Koblenz, Germany, October 23-26, 2023, pages 611–</i>	outcome-based feedback . <i>CoRR</i> , abs/2211.14275.	758
711	626. ACM.		
712	Minpeng Liao, Chengxi Li, Wei Luo, Jing Wu, and	Ziyu Wan, Xidong Feng, Muning Wen, Stephen Marcus	759
	Kai Fan. 2024. MARIO: math reasoning with code	McAlee, Ying Wen, Weinan Zhang, and Jun Wang.	760
	interpreter output - A reproducible pipeline . In <i>Find-</i>	2024. Alphazero-like tree-search can guide large	761
	<i>ings of the Association for Computational Linguistics,</i>	language model decoding and training . In <i>Forty-</i>	762
		<i>first International Conference on Machine Learning,</i>	763
		<i>ICML 2024, Vienna, Austria, July 21-27, 2024. Open-</i>	764
		<i>Review.net.</i>	765

766	Chaojie Wang, Yanchen Deng, Zhiyi Lv, Zeng Liang,	Li, Adrian Weller, and Weiyang Liu. 2024b. Meta-	822
767	Jujie He, Shuicheng Yan, and Bo An. 2024a. Q*:	math: Bootstrap your own mathematical questions	823
768	Improving multi-step reasoning for llms with deliber-	for large language models . In <i>The Twelfth Inter-</i>	824
769	ative planning . <i>CoRR</i> , abs/2406.14283.	<i>national Conference on Learning Representations,</i>	825
770	Peiyi Wang, Lei Li, Zhihong Shao, Runxin Xu, Damai	<i>ICLR 2024, Vienna, Austria, May 7-11, 2024</i> . Open-	826
771	Dai, Yifei Li, Deli Chen, Yu Wu, and Zhifang Sui.	Review.net.	827
772	2024b. Math-shepherd: Verify and reinforce llms	Eric Zelikman, Yuhuai Wu, Jesse Mu, and Noah D.	828
773	step-by-step without human annotations . In <i>Proceed-</i>	Goodman. 2022. Star: Bootstrapping reasoning with	829
774	<i>ings of the 62nd Annual Meeting of the Association</i>	reasoning . In <i>Advances in Neural Information Pro-</i>	830
775	<i>for Computational Linguistics (Volume 1: Long Pa-</i>	<i>cessing Systems 35: Annual Conference on Neural</i>	831
776	<i>pers)</i> , ACL 2024, Bangkok, Thailand, August 11-16,	<i>Information Processing Systems 2022, NeurIPS 2022,</i>	832
777	2024, pages 9426–9439. Association for Computa-	<i>New Orleans, LA, USA, November 28 - December 9,</i>	833
778	tional Linguistics.	2022.	834
779	Yu Wang, Nan Yang, Liang Wang, and Furu Wei.	Di Zhang, Jianbo Wu, Jingdi Lei, Tong Che, Jiatong	835
780	2025. Examining false positives under inference	Li, Tong Xie, Xiaoshui Huang, Shufei Zhang, Marco	836
781	scaling for mathematical reasoning. <i>arXiv preprint</i>	Pavone, Yuqiang Li, Wanli Ouyang, and Dongzhan	837
782	<i>arXiv:2502.06217</i> .	Zhou. 2024a. Llama-berry: Pairwise optimization	838
783	Jason Wei, Xuezhi Wang, Dale Schuurmans, Maarten	for o1-like olympiad-level mathematical reasoning .	839
784	Bosma, Brian Ichter, Fei Xia, Ed H. Chi, Quoc V. Le,	<i>CoRR</i> , abs/2410.02884.	840
785	and Denny Zhou. 2022. Chain-of-thought prompting	Hanning Zhang, Pengcheng Wang, Shizhe Diao, Yong	841
786	elicits reasoning in large language models . In <i>Ad-</i>	Lin, Rui Pan, Hanze Dong, Dylan Zhang, Pavlo	842
787	<i>vances in Neural Information Processing Systems 35:</i>	Molchanov, and Tong Zhang. 2024b. Entropy-	843
788	<i>Annual Conference on Neural Information Process-</i>	regularized process reward model . <i>CoRR</i> ,	844
789	<i>ing Systems 2022, NeurIPS 2022, New Orleans, LA,</i>	abs/2412.11006.	845
790	<i>USA, November 28 - December 9, 2022</i> .	Lunjun Zhang, Arian Hosseini, Hritik Bansal, Mehran	846
791	Yangzhen Wu, Zhiqing Sun, Shanda Li, Sean Welleck,	Kazemi, Aviral Kumar, and Rishabh Agarwal. 2024c.	847
792	and Yiming Yang. 2024. Inference scaling laws: An	Generative verifiers: Reward modeling as next-token	848
793	empirical analysis of compute-optimal inference for	prediction . <i>CoRR</i> , abs/2408.15240.	849
794	llm problem-solving. In <i>The 4th Workshop on Math-</i>	Zhenru Zhang, Chujie Zheng, Yangzhen Wu, Beichen	850
795	<i>ematical Reasoning and AI at NeurIPS'24</i> .	Zhang, Runji Lin, Bowen Yu, Dayiheng Liu, Jin-	851
796	Wei Xiong, Hanning Zhang, Nan Jiang, and Tong Zhang.	gren Zhou, and Junyang Lin. 2025. The lessons of	852
797	2024. An implementation of generative prm. https:	developing process reward models in mathematical	853
798	/github.com/RLHFlow/RLHF-Reward-Modeling .	reasoning. <i>arXiv preprint arXiv:2501.07301</i> .	854
799	An Yang, Baosong Yang, Beichen Zhang, Binyuan	Yuchen Zhuang, Xiang Chen, Tong Yu, Saayan Mitra,	855
800	Hui, Bo Zheng, Bowen Yu, Chengyuan Li, Dayi-	Victor S. Bursztyn, Ryan A. Rossi, Somdeb Sarkhel,	856
801	heng Liu, Fei Huang, Haoran Wei, Huan Lin, Jian	and Chao Zhang. 2024. Toolchain*: Efficient ac-	857
802	Yang, Jianhong Tu, Jianwei Zhang, Jianxin Yang,	tion space navigation in large language models with	858
803	Jiaxi Yang, Jingren Zhou, Junyang Lin, Kai Dang,	a* search . In <i>The Twelfth International Conference</i>	859
804	Keming Lu, Keqin Bao, Kexin Yang, Le Yu, Mei	<i>on Learning Representations, ICLR 2024, Vienna,</i>	860
805	Li, Mingfeng Xue, Pei Zhang, Qin Zhu, Rui Men,	<i>Austria, May 7-11, 2024</i> . OpenReview.net.	861
806	Runji Lin, Tianhao Li, Tingyu Xia, Xingzhang Ren,		
807	Xuancheng Ren, Yang Fan, Yang Su, Yichang Zhang,		
808	Yu Wan, Yuqiong Liu, Zeyu Cui, Zhenru Zhang, and		
809	Zihan Qiu. 2024. Qwen2.5 technical report . <i>CoRR</i> ,		
810	abs/2412.15115.		
811	Fei Yu, Anningzhe Gao, and Benyou Wang. 2024a.		
812	Ovm, outcome-supervised value models for planning		
813	in mathematical reasoning . In <i>Findings of the Asso-</i>		
814	<i>ciation for Computational Linguistics: NAACL 2024,</i>		
815	<i>Mexico City, Mexico, June 16-21, 2024</i> , pages 858–		
816	875. Association for Computational Linguistics.		
817	Fei Yu, Yingru Li, and Benyou Wang. 2025. Scaling		
818	flaws of verifier-guided search in mathematical rea-		
819	soning. <i>arXiv preprint arXiv:2502.00271</i> .		
820	Longhui Yu, Weisen Jiang, Han Shi, Jincheng Yu,		
821	Zhengying Liu, Yu Zhang, James T. Kwok, Zhenguo		

A Experiment Details

A.1 Baselines.

Outcome Reward Model (ORM, Cobbe et al., 2021). The vanilla ORM assigns a score to the entire solution as the final reward. We train ORMs through outcome supervision. Following (Cobbe et al., 2021), we replicate the binary correctness label $r_t \in \{0, 1\}$ across the entire solution sequence. The reward head is then trained to predict reward scores for each token, enhancing robustness.

Process Reward Model (PRM, Lightman et al., 2024; Uesato et al., 2022). The vanilla PRM assigns scores to each step along a solution path. For training stability, we place the reward label r_t at the last token of each step. In other words, for t -th step-level sequence, the label vector $\mathbf{y}_t = [0, 0, \dots, 0, r_t]$.

Math-shepherd PRM (Wang et al., 2024b). Different from the vanilla PRM, Math-Shepherd PRM uses Monte-Carlo Estimation to annotate step labels. This estimation is essentially considered as training a value model (Zhang et al., 2025). In our experiments, we first sample 15 solutions for each query. Then, for each intermediate step, we sample 8 rollouts. We provide a detailed description of this method in Section 4.2.

ER-PRM (Zhang et al., 2024b). Similar to Math-Shepherd PRM, ER-PRM integrates entropy-regularized step labels to train the supervision model. After Monte-Carlo sampling, ER-PRM calculates the label for the t -th step according to the following equation:

$$\text{label}(s_t) = \frac{1}{\eta} \ln \mathbb{E}_{\tau \sim \pi} e^{\eta y(\tau)}$$

where τ represents the complete rollout starting from the step s_t , π represents the LLM generator, and $y(\cdot)$ represents the final correctness of the solution τ .

A.2 BiRM Training Details

In the BiRM training phase, we collect problems from the original GSM8K and MATH dataset. Then we use LLM generators to sample 15 trajectories per query, resulting in a training set of approximately 225,000 solutions for each base model. We annotate reward and value labels using the Deepseek-V3 (DeepSeek-AI et al., 2024) and Math-shepherd soft-label methods, respectively.

We set training labels on the last token of each step, following (Wang et al., 2024b). The coefficient c in Equation 4 is set to 1.0.

A.3 Evaluation Metrics

At test-time, BiRM estimates both reward scores and value scores for partial solutions at the same time. We follow Equation 3 to calculate the final score. The coefficient β for different base models on GSM8K, MATH-500, and Gaokao2023 are set to $\beta_{\text{Qwen2.5-3B}} = \{1.0, 2.5, 2.0\}$, $\beta_{\text{Qwen2.5-7B}} = \{1.5, 3.0, 3.5\}$, $\beta_{\text{Llama3.1-8B}} = \{2.5, 1.0, 3.5\}$ respectively.

Best-of-N Sampling. For a given question q , we sample multiple rollouts from the LLM, resulting in a candidate set of N reasoning paths $\mathcal{T} = \{\tau_1, \tau_2, \dots, \tau_N\}$. Subsequently, an additional supervision model \mathcal{R} , such as PRM, is used to score each candidate path, yielding $\mathcal{R}(\tau_i)$, where $i \in \{1, 2, \dots, N\}$. The candidate with the highest score represents the most promising solution and is selected as the final output:

$$\tau^* = \arg \max_{\tau \in \{\tau_1, \tau_2, \dots, \tau_N\}} \mathcal{R}(\tau)$$

Beam Search. We present all search results from the main experiment in Table 5, Table 6, and Table 7. The procedure of the step-level beam search is as follows: We first set the total sampling size K and beam size b (K should be divisible by b). In each round, we only expand one step forward. For a given query, we sample K rollouts in the first round. Then, we use the supervision model \mathcal{M} to re-rank these candidates and select the top b rollouts for the next step. Starting from the second round, we expand $\frac{K}{b}$ trajectories for each candidate, getting K candidates in total. We repeat the re-ranking process until a final answer is found or the maximum step count is reached. A detailed pseudocode is provided in 1.

B Step Label Annotation Details

B.1 Dataset preprocessing

Before the SFT phase, we first preprocess the training data and restructure the delimiters at different levels of granularity. This is because we discover that original solution paths contain numerous meaningless text segments, which hinder the effective learning of process supervision models. Similar findings are reported by (Liao et al., 2024). To address this, we utilize Deepseek-V3 to clean the

MATH subset in the MetaMath dataset, reannotate the delimiters, and ensure that each step is logically complete and meaningful. The prompt template for data preprocessing is shown in Figure 5.

B.2 Reward Label Annotation

We also use Deepseek-V3 to annotate the correctness of each step (i.e., reward label) in our experiments. The prompt template is provided in Figure 6.

Algorithm 1 Step-Level Beam Search

```

1: Input: Question  $q$ , Total Sampling Size  $K$ , Beam size  $b$ ,
   Maximum step count  $T$ 
2: Output: Best solution path for  $q$ 
3: Model: Generator  $\pi$  and BiRM  $\mathcal{M}$ 
4: procedure STEPLEVELBEAMSEARCH( $q, K, b$ )
5:   Initialize partial solutions  $\mathbb{T} \leftarrow \{\}$ 
6:   Sample initial steps  $\{\tau_1^1, \tau_2^1, \dots, \tau_K^1\}$ 
7:   Estimate scores  $\{s_1^1, s_2^1, \dots, s_K^1\}$  for each step
8:   Select top  $b$  scored steps and add to  $\mathbb{T}$ 
9:    $t \leftarrow 1$ 
10:  while solutions in  $\mathbb{T}$  are not complete and  $t < T$  do
11:    New candidate solutions  $\mathbb{T}_{\text{new}} \leftarrow \{\}$ 
12:    Scores  $\mathcal{S} \leftarrow \{\}$ 
13:    for each partial solution  $\tau^{[1:t]}$  in  $\mathbb{T}$  do
14:      for  $i = 1$  to  $K/b$  do
15:         $\tau_i^{[1:t+1]} \sim \pi(\tau^{[1:t]}, q)$ 
16:         $s_i^{[1:t+1]} = \mathcal{M}(\tau_i^{[1:t+1]}, q)$ 
17:         $\mathbb{T}_{\text{new}} \leftarrow \mathbb{T}_{\text{new}} + \tau_i^{[1:t+1]}$ 
18:         $\mathcal{S} \leftarrow \mathcal{S} + s_i^{[1:t+1]}$ 
19:      end for
20:    end for
21:     $\mathbb{T}_{\text{new}} \leftarrow$  top  $b$  scored partial solutions in  $\mathbb{T}_{\text{new}}$ 
22:     $\mathbb{T} \leftarrow \mathbb{T}_{\text{new}}$ 
23:     $t \leftarrow t + 1$ 
24:  end while
25:  return solution with the highest score in  $\mathbb{T}$ 
26: end procedure

```

Total Size	Beam Size	GSM8K			MATH-500			Gaokao2023		
		OVM	PRM	BiRM	OVM	PRM	BiRM	OVM	PRM	BiRM
$K = 4$	4	81.50 \pm 0.45	82.11 \pm 0.28	81.96 \pm 0.38	48.60 \pm 0.16	48.13 \pm 1.20	47.27 \pm 0.96	33.85 \pm 0.53	33.16 \pm 0.74	33.07 \pm 0.44
	2	82.97 \pm 0.14	81.53 \pm 0.72	82.76 \pm 0.28	48.27 \pm 1.61	49.27 \pm 0.81	50.07 \pm 0.84	35.15 \pm 1.80	34.55 \pm 0.56	36.10 \pm 0.76
	1	80.82 \pm 0.51	80.57 \pm 0.77	81.93 \pm 0.64	47.60 \pm 0.59	47.80 \pm 1.14	47.67 \pm 0.57	35.58 \pm 1.29	34.89 \pm 1.56	32.47 \pm 0.56
$K = 8$	8	83.70 \pm 0.73	83.65 \pm 0.25	84.41 \pm 0.35	49.33 \pm 0.38	50.13 \pm 1.00	50.07 \pm 0.66	35.93 \pm 1.17	34.63 \pm 0.49	35.06 \pm 1.39
	4	84.61 \pm 0.56	83.93 \pm 0.16	85.11 \pm 0.40	48.87 \pm 0.68	50.87 \pm 0.41	52.53 \pm 0.90	36.10 \pm 1.10	37.92 \pm 0.76	37.92 \pm 0.21
	2	84.10 \pm 0.36	83.17 \pm 0.25	84.00 \pm 0.39	50.07 \pm 1.06	50.33 \pm 0.94	50.67 \pm 1.64	35.32 \pm 1.10	37.58 \pm 1.09	36.97 \pm 1.07
	1	83.27 \pm 0.40	82.66 \pm 1.05	82.99 \pm 0.23	48.47 \pm 2.03	49.67 \pm 0.25	49.93 \pm 0.41	33.42 \pm 1.59	35.24 \pm 0.32	35.06 \pm 1.12
$K = 20$	20	85.27 \pm 0.04	85.65 \pm 0.50	86.13 \pm 0.11	52.13 \pm 1.15	53.20 \pm 0.59	53.33 \pm 1.48	36.54 \pm 0.44	35.67 \pm 1.56	36.10 \pm 1.85
	10	86.73 \pm 0.65	84.66 \pm 0.34	86.91 \pm 0.25	53.00 \pm 0.16	54.27 \pm 0.77	55.00 \pm 0.65	37.66 \pm 1.48	38.35 \pm 1.24	37.23 \pm 1.38
	5	86.23 \pm 0.28	84.86 \pm 0.54	86.28 \pm 0.33	52.20 \pm 0.59	53.40 \pm 0.85	54.27 \pm 0.52	36.88 \pm 1.06	37.49 \pm 0.86	37.58 \pm 0.24
	4	86.20 \pm 0.16	84.76 \pm 0.22	85.04 \pm 0.19	51.73 \pm 0.66	51.80 \pm 0.75	53.60 \pm 0.75	37.49 \pm 0.88	35.41 \pm 1.00	39.05 \pm 1.17
	2	85.32 \pm 0.25	84.74 \pm 0.36	85.19 \pm 0.19	49.00 \pm 0.33	50.33 \pm 1.00	51.80 \pm 0.49	35.67 \pm 0.44	35.84 \pm 0.97	36.62 \pm 0.85
	1	83.60 \pm 0.09	82.56 \pm 0.74	84.23 \pm 0.39	49.00 \pm 0.91	50.67 \pm 0.90	50.87 \pm 0.84	34.29 \pm 1.29	34.46 \pm 1.41	37.06 \pm 2.51
$K = 100$	50	87.29 \pm 0.22	85.87 \pm 0.64	87.34 \pm 0.22	52.87 \pm 0.82	53.87 \pm 0.19	55.13 \pm 0.34	37.06 \pm 0.74	37.40 \pm 0.97	38.96 \pm 0.92
	25	87.54 \pm 0.26	85.52 \pm 0.80	87.64 \pm 0.65	53.00 \pm 1.50	53.20 \pm 0.33	54.73 \pm 0.75	38.10 \pm 1.17	37.75 \pm 1.09	38.18 \pm 0.21
	10	85.90 \pm 0.33	84.51 \pm 0.77	86.71 \pm 0.37	51.27 \pm 1.32	49.80 \pm 0.57	53.40 \pm 1.23	38.01 \pm 1.05	37.92 \pm 1.10	37.40 \pm 1.85

Table 5: Qwen2.5-3B performance of beam search on GSM8K, MATH-500 and Gaokao2023.

Total Size	Beam Size	GSM8K			MATH-500			Gaokao2023		
		OVM	PRM	BiRM	OVM	PRM	BiRM	OVM	PRM	BiRM
$K = 4$	4	86.10 \pm 0.52	86.48 \pm 0.43	87.04 \pm 0.16	55.73 \pm 1.52	53.73 \pm 1.51	57.13 \pm 1.15	40.09 \pm 0.12	41.13 \pm 1.24	43.90 \pm 0.92
	2	86.20 \pm 0.53	86.00 \pm 0.25	86.99 \pm 0.13	55.53 \pm 0.47	55.80 \pm 1.72	55.87 \pm 0.93	42.77 \pm 0.24	42.68 \pm 0.88	43.55 \pm 0.61
	1	85.65 \pm 1.01	85.04 \pm 0.50	86.76 \pm 0.31	53.80 \pm 1.77	54.93 \pm 1.46	56.33 \pm 1.04	41.47 \pm 0.74	43.98 \pm 0.12	44.50 \pm 1.17
$K = 8$	8	86.73 \pm 0.62	88.12 \pm 0.36	88.93 \pm 0.49	58.27 \pm 0.50	58.20 \pm 0.71	58.13 \pm 0.90	44.19 \pm 0.76	44.24 \pm 0.68	45.11 \pm 0.68
	4	88.63 \pm 0.19	87.89 \pm 0.73	89.36 \pm 0.22	57.60 \pm 0.85	59.07 \pm 1.32	59.53 \pm 1.24	43.72 \pm 0.86	45.63 \pm 1.21	46.84 \pm 0.65
	2	88.55 \pm 0.37	87.45 \pm 0.19	88.30 \pm 0.53	57.00 \pm 0.43	57.20 \pm 1.56	58.67 \pm 1.27	43.98 \pm 1.41	44.68 \pm 0.56	45.45 \pm 0.97
	1	87.57 \pm 0.38	86.45 \pm 0.40	87.47 \pm 0.19	54.67 \pm 1.09	57.27 \pm 1.23	57.73 \pm 1.32	44.24 \pm 1.09	44.94 \pm 1.27	43.38 \pm 0.52
$K = 20$	20	86.33 \pm 0.38	88.65 \pm 0.09	90.04 \pm 0.58	59.07 \pm 0.82	59.60 \pm 0.59	60.33 \pm 0.68	44.76 \pm 1.38	45.89 \pm 1.21	47.71 \pm 0.44
	10	90.40 \pm 0.18	89.18 \pm 0.42	90.40 \pm 0.65	58.73 \pm 1.16	61.53 \pm 1.05	62.27 \pm 1.09	45.19 \pm 0.76	48.14 \pm 0.74	48.23 \pm 0.12
	5	90.30 \pm 0.12	88.98 \pm 0.26	90.60 \pm 0.28	57.53 \pm 0.34	59.40 \pm 0.49	60.73 \pm 0.19	45.45 \pm 0.64	47.36 \pm 1.17	47.36 \pm 1.22
	4	89.56 \pm 0.25	87.52 \pm 0.62	89.94 \pm 0.09	56.53 \pm 0.90	58.40 \pm 1.28	59.67 \pm 0.34	45.02 \pm 0.53	45.54 \pm 0.88	48.31 \pm 2.21
	2	88.55 \pm 0.06	88.05 \pm 0.07	89.69 \pm 0.34	56.93 \pm 0.90	57.47 \pm 1.11	58.87 \pm 0.62	43.72 \pm 1.07	44.59 \pm 0.74	47.62 \pm 0.96
	1	87.79 \pm 0.22	86.96 \pm 0.47	88.07 \pm 0.38	56.27 \pm 0.34	57.73 \pm 1.52	58.33 \pm 0.77	42.68 \pm 1.71	45.63 \pm 0.74	45.80 \pm 0.86
$K = 100$	50	91.00 \pm 0.22	88.32 \pm 0.57	91.28 \pm 0.12	60.13 \pm 0.47	60.73 \pm 0.34	62.53 \pm 0.77	46.84 \pm 0.44	48.31 \pm 0.85	49.96 \pm 0.32
	25	91.18 \pm 0.53	88.40 \pm 0.21	91.66 \pm 0.33	58.40 \pm 1.23	59.27 \pm 0.34	62.00 \pm 0.98	46.32 \pm 0.53	47.97 \pm 0.12	47.62 \pm 0.68
	10	89.97 \pm 0.25	88.15 \pm 0.16	91.00 \pm 0.09	57.47 \pm 1.32	59.20 \pm 1.31	61.20 \pm 0.43	43.64 \pm 0.42	46.93 \pm 0.49	49.00 \pm 0.32

Table 6: Qwen2.5-7B performance of beam search on GSM8K, MATH-500 and Gaokao2023.

Total Size	Beam Size	GSM8K			MATH-500			Gaokao2023		
		OVM	PRM	BiRM	OVM	PRM	BiRM	OVM	PRM	BiRM
$K = 4$	4	71.44 \pm 0.36	71.65 \pm 0.33	71.37 \pm 0.41	37.53 \pm 0.66	38.67 \pm 0.50	38.07 \pm 0.68	23.81 \pm 0.65	24.94 \pm 0.21	23.29 \pm 0.86
	2	72.76 \pm 0.41	71.11 \pm 1.11	72.91 \pm 0.80	38.53 \pm 2.22	39.87 \pm 0.68	40.73 \pm 0.52	23.90 \pm 1.10	25.11 \pm 1.36	26.06 \pm 0.68
	1	70.74 \pm 0.16	68.99 \pm 0.67	71.57 \pm 0.38	36.80 \pm 0.75	39.60 \pm 0.85	39.33 \pm 0.98	23.20 \pm 0.74	24.33 \pm 0.74	24.59 \pm 0.65
$K = 8$	8	76.52 \pm 0.45	75.92 \pm 0.20	76.90 \pm 0.53	39.93 \pm 1.48	39.00 \pm 0.59	41.27 \pm 0.09	25.63 \pm 1.07	25.71 \pm 0.56	26.75 \pm 0.56
	4	77.36 \pm 0.47	75.84 \pm 0.54	78.32 \pm 0.55	40.20 \pm 0.49	40.13 \pm 1.64	43.27 \pm 0.57	25.19 \pm 1.29	26.49 \pm 1.53	27.45 \pm 1.56
	2	75.51 \pm 0.49	73.79 \pm 1.02	76.17 \pm 0.04	39.07 \pm 0.50	39.80 \pm 1.85	41.47 \pm 1.23	25.02 \pm 0.96	26.58 \pm 2.04	26.23 \pm 1.18
	1	74.00 \pm 0.66	72.40 \pm 0.57	74.60 \pm 0.62	37.27 \pm 1.64	40.13 \pm 1.57	41.47 \pm 0.77	23.72 \pm 0.74	25.28 \pm 1.59	25.80 \pm 1.17
$K = 20$	20	79.93 \pm 0.22	79.23 \pm 0.48	80.46 \pm 0.29	41.53 \pm 0.84	41.00 \pm 0.75	44.13 \pm 0.19	26.15 \pm 0.74	25.71 \pm 0.21	27.62 \pm 0.44
	10	81.40 \pm 0.25	78.82 \pm 0.22	81.73 \pm 0.62	40.73 \pm 0.68	41.60 \pm 0.86	44.27 \pm 0.34	25.97 \pm 0.97	28.57 \pm 1.18	29.18 \pm 0.86
	5	79.76 \pm 0.37	76.90 \pm 0.35	81.00 \pm 0.20	40.80 \pm 0.71	42.07 \pm 1.09	43.93 \pm 1.00	27.01 \pm 0.37	28.14 \pm 0.98	28.57 \pm 0.73
	4	79.56 \pm 0.26	76.02 \pm 0.36	80.16 \pm 0.64	40.13 \pm 1.55	42.00 \pm 1.07	44.00 \pm 0.43	24.33 \pm 0.32	28.23 \pm 0.12	26.93 \pm 0.80
	2	77.96 \pm 0.60	75.39 \pm 1.04	79.53 \pm 1.08	39.07 \pm 0.90	39.53 \pm 0.50	42.40 \pm 0.65	26.58 \pm 1.41	26.75 \pm 0.85	26.84 \pm 0.12
	1	76.27 \pm 0.80	73.24 \pm 1.02	78.17 \pm 0.87	39.27 \pm 0.82	40.00 \pm 1.82	41.53 \pm 1.36	25.54 \pm 0.24	27.36 \pm 1.44	27.27 \pm 0.56
$K = 100$	50	82.71 \pm 0.11	80.34 \pm 0.84	85.39 \pm 0.52	41.07 \pm 0.50	42.33 \pm 0.66	46.13 \pm 0.98	26.23 \pm 2.02	29.61 \pm 0.37	30.65 \pm 0.37
	25	82.71 \pm 0.61	78.44 \pm 0.43	84.53 \pm 0.60	40.93 \pm 1.32	42.00 \pm 0.71	44.07 \pm 0.68	25.37 \pm 0.61	28.14 \pm 0.86	29.00 \pm 0.61
	10	81.10 \pm 0.46	77.81 \pm 0.93	83.35 \pm 0.70	39.87 \pm 0.77	40.27 \pm 1.06	45.00 \pm 0.16	25.80 \pm 0.32	27.19 \pm 0.86	29.70 \pm 0.80

Table 7: Llama3.1-8B performance of beam search on GSM8K, MATH-500 and Gaokao2023.

SFT Dataset Preprocessing

You are an expert math examiner, skilled at transforming complex mathematical solution steps into clear formats. Your task is to insert the symbol `<step_end>` to mark the end of each step in the following math problem's solution. A step should represent a complete statement, structure or calculation process. You must not omit any original content and only insert this symbol at the end of each step. Your output should only include the revised solution, without any additional text.

Figure 5: The prompt template for MetaMath dataset preprocessing.

Reward Label Annotation

You are an expert math examiner. Your task is to review the student's solution and evaluate each step. Mark a step as correct only if it is based on accurate premises and contributes to solving the problem. Mark it as unnecessary if it is logically valid but does not aid in solving the problem. Your judgments should include a very concise analysis of each step and the final judgement.

You must provide your evaluations in JSON format like:

```
{"step_1": {"analysis": "<concise analysis of the step>", "judgement":  
"<correct/incorrect/unnecessary>"}, "step_2": {...}, ...}
```

Below is the question, reference answer, and student's solution that you need to evaluate. Note that the student's solution does not need to match the reference solution exactly.

[Question]
{question}

[Reference Answer]
{answer}

[Solution]
{solution}

Now, provide your evaluations in JSON format.

Figure 6: The prompt template for reward label annotation.